

Convolutional Neural Network for Accurate Crowd Counting and Density Estimation

by **Saeed Amirgholipour Kasmani**

Thesis submitted in fulfilment of the requirements for
the degree of

Doctor of Philosophy

under the supervision of Professor Xiangjian (Sean) He

University of Technology Sydney
Faculty of Engineering and IT

February 2021

Certificate of Original Authorship Template

Graduate research students are required to make a declaration of original authorship when they submit the thesis for examination and in the final bound copies. Please note, the Research Training Program (RTP) statement is for all students. The Certificate of Original Authorship must be placed within the thesis, immediately after the thesis title page.

Required wording for the certificate of original authorship

CERTIFICATE OF ORIGINAL AUTHORSHIP

I, **Saeed Amirgholipour Kasmani** declare that this thesis, is submitted in fulfilment of the requirements for the award of **PhD**, in the **SEDE/FEIT** at the University of Technology Sydney.

This thesis is wholly my own work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

**If applicable, the above statement must be replaced with the collaborative doctoral degree statement (see below).*

**If applicable, the Indigenous Cultural and Intellectual Property (ICIP) statement must be added (see below).*

This research is supported by the Australian Government Research Training Program.

Signature: Production Note:
 Signature removed prior to publication.

Date: 26/02/2021

Collaborative doctoral research degree statement

I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of the requirements for a degree at any other academic institution except as fully acknowledged within the text. This thesis is the result of a Collaborative Doctoral Research Degree program with **CSIRO\Data 61**.

Indigenous Cultural and Intellectual Property (ICIP) statement

This thesis includes Indigenous Cultural and Intellectual Property (ICIP) belonging to *[insert relevant language, tribal or nation group(s) or communities]*, custodians or traditional owners. Where I have used ICIP, I have followed the relevant protocols and consulted with appropriate Indigenous people/communities about its inclusion in my thesis. ICIP rights are Indigenous heritage and will always remain with these groups. To use, adapt or reference the ICIP contained in this work, you will need to consult with the relevant Indigenous groups and follow cultural protocols.

ABSTRACT

Convolutional Neural Network for Accurate Crowd Counting and Density Estimation

by

Saeed Amirgholipour Kasmani

Nowadays, crowd and object counting has become an important task for a variety of applications, such as traffic control, public safety, urban planning, and video surveillance. It has also become a crucial part of building a high-level monitoring system such as video surveillance and crowd analysis. In these cases, dynamic crowd monitoring and analysis is extremely important for control management and social safety.

Like the other computer vision issues, crowd counting and density estimation come with various kinds of challenges such as high clutters, occlusions, non-uniform distributions of objects or people, and intra-scene and inter-scene variations in appearance. Researchers and industrial partners have attempted to design and develop many sophisticated models to address various issues that exist in crowd counting. Especially in recent years, the number of researches in the crowd counting era became overwhelming with the domination of deep-learning and Convolution Neural Networks (CNNs) based models in various computer vision tasks. In this thesis, we revisit the crowd counting and propose various novel solutions to this problem.

At first, we propose an Adaptive Counting Convolutional Neural Network (A-CCNN) and consider the scale variation of objects in a frame adaptively to improve the accuracy of counting. Our method takes advantages of contextual information to provide more accurate and adaptive density maps and crowd counting in a scene. Then, we focus on CNN pruning to further enhance the crowd counting models for real-time application and increase the performance of CCNN model. Thus, a new

pruning strategy is proposed by considering the contributions of various filters to the final result. The filters in the original CCNN model are grouped into positive, negative, and irrelevant types. We prune the irrelevant filters, of which feature maps contain little information, and the negative filters determined by a mask learned from the training dataset. Our solution improves the results of the counting model without fine-tuning or retraining the pruned model. Finally, we propose a novel Pyramid Density-Aware Attention-based network, abbreviated as PDANet, which leverages the attention, pyramid scale feature and two branch decoder modules for density-aware crowd counting. The PDANet utilises these modules to extract different scale features, focus on the relevant information, and suppress the misleading ones. Extensive evaluations conducted on the challenging benchmark datasets well demonstrate the superior performance of the proposed models in terms of the accuracy of counting as well as generated density maps over the well-known state-of-the-art approaches.

Dissertation directed by Professor Xiangjian (Sean) He
School of Electrical and Data Engineering

Dedication

I owe thanks to a very special person, my wife, Mozhgan for her continued and unfailing love, support and understanding during my pursuit of Ph.D degree that made the completion of thesis possible. This thesis work is dedicated to you who were always around at times I thought that it was impossible to continue, and you helped me to keep things in perspective. I greatly value her contribution and deeply appreciate her belief in me.

This work is also dedicated to my mother, Soraya, who had always loved me unconditionally and whose good examples have taught me to work hard for the things that I aspire to achieve.

Acknowledgements

At this moment of accomplishment, I would like to express my deep appreciation to my wonderful supervisor, Prof Sean He, who has been guiding, supporting, and encouraging me throughout my PhD candidature journey. He gave me care and love like a father, which I can never forget. I am also greatly indebted to my co-supervisor, Dr Wenjing Jia, for her involvement and guidance throughout my PhD. She offered me her positive appreciations, constructive criticisms, and valuable advices during my PhD research which helped me tackle the challenges and complete the research work. I extend thanks to Prof Dadong Wang, my co-supervisor from CSIRO/DATA 61, for his time, encouragement, and expertise throughout this journey. I was able to accomplish this research due to the tremendous and unforgettable supports from my supervisory panel.

I need to express my deep gratitude to UTS and CSIRO/ Data61 to sponsor support and scholarships. I am thankful to the University of Technology Sydney, to give me permission to conduct PhD studies and provide me with state-of-the-art facilities to do my experiments. I am thankful to all staff at the School of Electrical and Data Engineering, FEIT, and GRS staffs of UTS for their kind assistance and support. It is my fortune to acknowledge my friends' support throughout the research tenure gratefully.

In the end, my wife, my family, and my parents, I sincerely thank you for believing me and showing faith in me. I would not have been able to focus on my research without the ongoing support and unconditional love of my wife.

Saeed Amirgholipour Kasmani
Sydney, Australia, 2021.

List of Publications

Journal Papers

- J-1. **Amirgholipour**, S., Jia, W., Liu, L., Fan, X., Wang, D. and He, X., 2021. PDANet: Pyramid Density-aware Attention Net for Accurate Crowd Counting. *Neurocomputing*, 451, pp.215-230.
- J-2. Liu, L., Jia, W., Jiang, J., **Amirgholipour**, S., Wang, Y., Zeibots, M. and He, X., 2020. Denet: A universal network for counting crowd with varying densities and scales. *IEEE Transactions on Multimedia*, 23, pp.1060-1068.
- J-3. Liu, L., **Amirgholipour**, S., Jiang, J., Jia, W., Zeibots, M. and He, X., 2019. Performance-enhancing network pruning for crowd counting. *Neurocomputing*, 360, pp.246-253.

Conference Papers

- C-1. **Amirgholipour**, S., He, X., Jia, W., Wang, D. and Zeibots, M., 2018, October. A-CCNN: Adaptive CCNN for density estimation and crowd counting. In 2018 25th IEEE International Conference on Image Processing (ICIP) (pp. 948-952). IEEE.

Contents

Certificate	ii
Abstract	iii
Dedication	v
Acknowledgments	vi
List of Publications	vii
List of Figures	xiii
Abbreviation	xvii
Notation	xviii
1 Introduction	1
1.1 Background	1
1.1.1 Related Works and Scope	2
1.2 Challenges	3
1.3 Research Specification	6
1.3.1 Research Stakeholders and Objectives	6
1.3.2 Research Significance	7
1.3.3 Ethics and Risk Consideration	8
1.4 Contributions	9
1.4.1 A-CCNN: Adaptive CCNN for Density Estimation and Crowd Counting [2]	9

1.4.2	Performance-Enhancing Network Pruning for Crowd Counting [82]	10
1.4.3	PDANet: Pyramid Density-aware Attention Net for Accurate Crowd Counting	10
1.5	Thesis Organization	11
2	Literature Review	13
2.1	Introduction	14
2.1.1	Organisation of This Chapter	14
2.2	Key Characteristic of Crowd Counting Solutions	16
2.2.1	Model Supervision Method	18
2.2.2	Model Training and Learning Procedure	19
2.2.3	Processing Manner	21
2.2.4	Crowd Counting Network Architectures	22
2.3	Datasets	28
2.3.1	Popular Datasets	29
2.3.2	Recently Proposed Datasets	31
2.3.3	Special Crowd Counting Datasets	32
2.3.4	Representative Object Counting Datasets in Other Fields	34
2.4	Evaluation Metrics	35
2.5	Evaluation and Discussion	36
2.5.1	Evaluation Results	37
2.5.2	Key Features of the Best Solutions	46
3	A-CCNN: Adaptive CCNN for Density Estimation and Crowd Counting	51

3.1	Introduction	51
3.1.1	The CCNN Architecture	53
3.2	Adaptive CCNN	54
3.2.1	Head Detection	55
3.2.2	Adaptive HP Selection by FIS	56
3.2.3	Training Parameters	57
3.3	Experimental Results	58
3.3.1	The UCSD Dataset	58
3.3.2	The UCF-CC Dataset	60
3.3.3	The Sydney Train Dataset	60
3.4	Conclusion	61
4	Performance-Enhancing Network Pruning for Crowd Counting	62
4.1	Introduction	62
4.2	Related Works	64
4.3	Network Pruning	66
4.3.1	Determining the Types of Filters	66
4.3.2	Pruning Filters and Feature Maps	70
4.3.3	Pruning of Different Layers	71
4.4	Experiments	72
4.4.1	Comparison with the Original CCNN Model	72
4.4.2	Comparison with Other Pruning Algorithms	74
4.4.3	Pruning Results on Other Crowd Counting Models	76
4.4.4	Impact of η	77

4.5 Conclusion	78
5 PDANet: Pyramid Density-aware Attention Net for Accurate Crowd Counting	79
5.1 Introduction	79
5.2 Related Works	84
5.3 Pyramid Density-aware Attention Net	85
5.3.1 Overview	86
5.3.2 Channel and Spatial based Attention Modules	88
5.3.3 Pyramid Feature Extractor (PFE)	90
5.3.4 Classification Module	92
5.3.5 Density Aware Decoder (DAD)	93
5.3.6 Implementation Details	96
5.3.7 Regression Loss and Ground Truth	96
5.4 Experiments	100
5.4.1 Data Augmentation	100
5.4.2 Experimental Results on the Shanghai Tech Dataset	100
5.4.3 Experimental Results on the WorldExpo10 Dataset	103
5.4.4 Experimental Results on the UCF Dataset	104
5.4.5 Experimental Results on the UCSD Dataset	106
5.5 Ablation Study	106
5.5.1 Density Map Visualization	107
5.5.2 Effectiveness of the PFE Module	110
5.5.3 Effectiveness of the Attention Module	114
5.5.4 Effectiveness of the Classification and DAD Modules	114

5.6 Conclusion	117
6 Conclusions and Future Work	118
6.1 Summary of the Thesis	119
6.2 Future Research	121
Bibliography	123

List of Figures

1.1	The existing tough challenges in the crowd counting task [35].	4
2.1	High level Classification of Crowd Counting solutions [35].	15
2.2	Three different architectures for existing crowd counting approaches [35].	23
3.1	The overview of our proposed A-CCNN crowd counting method. For an input image, our A-CCNN first estimates head size and corresponding position, and then utilises a fuzzy engine to apply separate CCNN models to each section to estimate the overall count.	53
3.2	The fuzzy inference engine, where the head size and corresponding position are the two inputs and the level of HPs for CCNN is the output.	56
4.1	Example of the feature maps in layer conv4. (a) Input image. (b) Filters activating mostly on targets. (c) Filters activating mostly on background. (d) Filters with nearly no activation.	67
4.2	Learning the mask from an annotated training image.	68
4.3	Model pruning with one mask.	69

4.4	The MAE results of the estimations obtained on different subsets of the UCF dataset [55] by pruning different layers with different ratios η	72
4.5	Examples of the density heat maps obtained with the original CCNN approach and our pruned CCNN model, where ground truth counts and estimation counts are shown underneath the images. (a) Input crowd images. (b) Density heat map obtained with the original CCNN. (c) Density heat map obtained with our pruned CCNN.	73
5.1	Examples of crowded and sparse images. (a) and (c) show an example of a highly crowded scene and a less crowded scene, respectively, while (b) and (d) show their corresponding congested areas.	81
5.2	The overview of our proposed PDANet network. This architecture contains a VGG16 based feature extractor, a Pyramid module, an Attention module, a Classification module, and a Decoder module.	86
5.3	Illustration of the attention module of our model. The top branch generates channel-based attention, while the bottom branch generates the spatial attention map.	88
5.4	The overview of the Pyramid Feature Extractor (PFE) module. The PFE module uses 1×1 and 3×3 dilated kernel convolutions with the GAP to extract features of different scales from the VGG16 features.	90
5.5	Illustration of the classification module of PDANet. It uses the global average pooling with a fully connected layer to determine the dense level of input scene.	93

5.6	The illustration of the DAD module. The input feature maps are fed to the two shared layers and then we use the two branches with three convolution layers to handle the dense and sparse areas within the scene.	94
5.7	Results of the estimated density maps of images from the Shanghai Tech part A dataset. We illustrate three test images (a0, b0, c0), their actual ground truth (a1, b1, c1), our estimated overall density maps (a2, b2, c2), our estimated density maps for dense areas (a3, b3, c3), and our estimated density maps for sparse areas and their crowd counts (a4, b4, c4).	108
5.8	Results of the estimated density maps of images from the Shanghai Tech part B dataset. We illustrate three test images (a0, b0, c0), their actual ground truth (a1, b1, c1), our estimated overall density maps (a2, b2, c2), our estimated density maps for dense areas (a3, b3, c3), and our estimated density maps for sparse areas and their crowd counts (a4, b4, c4).	110
5.9	Results of the estimated density maps of images from the UCF CC 50 dataset [55]. We present three test images (a0, b0, c0), their actual ground truth (a1, b1, c1), our estimated overall density maps (a2, b2, c2), our estimated density maps for dense areas (a3, b3, c3), and our estimated density maps for sparse areas and their crowd counts (a4, b4, c4).	111
5.10	Comparison of MAE and MSE results between various numbers of GAP layers on the UCF CC 50 crowd counting [55].	112

- 5.11 Comparison of MAE and MSE results between various numbers of GAP layers on the Shanghai Tech partA [183]. . 112
- 5.12 Comparison of MAE and MSE results between various numbers of GAP layers on the Shanghai Tech part B [183]. . 113

Abbreviation

CNN - Convolutional Neural Network

MAE - Mean Absolute Error

RMSE - Root Means Square error

Weak-Sup: Weak Supervised

Self-Sup: Self Supervised

Semi-Sup: Semi Supervised

Un-Sup: Un Supervised.

Fully-Sup: Fully Supervised

Fib: Full Image-Based

Pb: Patch-based

DOF: Degree Of Freedom

HPs: Hyper Parameters

Nomenclature and Notation

Capital letters denote matrices.

Lower-case alphabets denote column vectors.

$(\cdot)^T$ denotes the transpose operation.

I_n is the identity matrix of dimension $n \times n$.

0_n is the zero matrix of dimension $n \times n$.

\mathbb{R} , \mathbb{R}^+ denote the field of real numbers, and the set of positive reals, respectively.